

**TITLE: METHOD OF PROVIDING DUPLICATE ORIGINAL FILE COPIES OF A SEARCHED TOPIC FROM MULTIPLE FILE TYPES DERIVED FROM THE WEB**

**FIELD OF THE INVENTION:**

The present disclosure involves methods for developing full text searches for searching multiple file types which are downloaded from the Web.

CROSS-REFERENCES TO RELATED APPLICATIONS:

This application is related to a co-pending application, USSN \_\_\_\_\_ entitled "Method For Searching Multiple File Types on a CD-ROM", which is  
5 incorporated herein by reference.

BACKGROUND OF THE INVENTION:

In present day commercial situations, many digital development software and computer companies work to deliver documentation to their customers in a number of different formats. These formats may show up in a  
10 number of different varieties, that is to say the document format may be on paper, for example, or Adobe Acrobat Portable Document Format (PDF) files, or Windows Help files, or Hypertext Markup Language (HTML) and also  
15 HTML help files.

The documentation provided to receivers, such as customers, is distributed and made available on, for example, paper documents, on CD ROMs, and on Web Servers.

Of course, it is desirable for a recipient or  
20 user to make a full text search of the received documents. However, users cannot perform full-text searches on paper documents, except through long, laborious reading and surveys of the documents. There is, however, software designated as "search engines" that  
25 exist in digital technology in order to search files that are distributed to users who download from the Web.

However, these search engines are limited in a number of ways in providing search capability when the document or received Web files involve multiple file  
30 types. Most of the existing search engines are designed only to search files of one particular format.

In this type of situation, then it would be necessary to convert all files in the Web documents or Web-received files into a common format. This common format would be the format which was compatible with the particular search engine available.

However, when files are converted into a format different from that in which they were originally created, much of the functionality for searching the original file is lost, and this includes navigating through the file and finding certain special graphics or other content in the file.

There are other types of search engines which are capable in a certain limited way of including search operations for multiple file types in the Web received file documentation. However, these search engines are unable to open all the file types at locations where the search terms appear and then be capable of moving from one such location to the next location within the document.

Thus, these other types of search engines require that the user first search with one particularly favorite engine and then refine the search using another search engine designed for the file type.

One example of a standard (not a full-text) search is what one can do in a product program such as Word. The operator tells Word to find a text string. Then Word starts reading the text in the document by reading each word one at a time beginning at a specified location and comparing the text against the string that was entered. Now, when Word finds a "hit" (match), then Word highlights the text and stops searching. If the operator chooses "Find Next" option, then the Word program repeats the process and continues the search beginning just past the current hit. However, this is considered pretty much of a brute force and very slow process of operation.

The present invention provides for the use of an existing search engine that is designed to support the searching of one particular file format (PDF, or Adobe® Acrobat® files). This can then be extended to allow the  
30 searching of virtually any other type of file format such as HTML, HTML Help, or Windows Help. The method and system accomplishes this by creating a PDF file "duplicate" consisting of the text from the file that the  
35 operator wants to search in order to allow the search engine to find the text in the duplicate that was created. Here then there is provided a link from each

awk\appl\041503L.doc

**SUMMARY OF THE INVENTION:**

The described method involves the handling of multiple files downloaded from the web which files may exist in quite different word formats which are not readily searchable for desired topics or word matches.

The present method and system involves a technique that converts the downloaded file types into a Portable Document Format which uses an Adobe Acrobat program to search Portable Document Format (PDF) files that contain the text extracted from files residing in other formats such as Windows Help, Hypertext Markup Language (HTML) Help, and HTML.

On each page of the PDF file there are hyperlinks that the user can select to open the original file at the corresponding location.

The method enables the user to search the collection of PDF files, including both files that were created as PDF files as well as the PDF files created from the text extracted from the files of other formats. The method uses the search engine from Verity that is distributed by Adobe® in order to search the Adobe® Acrobat® portable document format files (PDF) which were downloaded from the Web. If the search targets include files of formats other than PDF, then the user is presented with pages within the PDF copy of the file in which the target text appears.

The user can navigate within the PDF copy using the "next hit" and "previous hit" program options. The text is visible to the user and is sufficient to help the user determine whether it is necessary or helpful to access the original file.

Each page of the PDF file carries a "button" that, when selected, opens the document in the original format at the location corresponding to the location displayed in the PDF copy. Both the PDF copy and the

5                   The indicated method includes software which is  
used to extract the text from Windows Help, HTML, and  
HTML Help files, and then create from that text the new  
files that can be converted by the standard Adobe  
software into PDF files with corresponding explanatory  
0 messages and buttons on every page in order to support  
the linking into the corresponding locations within the  
original files.

15 corresponding locations within the original files.

**BRIEF DESCRIPTION OF THE DRAWINGS:**

Fig. 1A is a block diagram illustrating the environmental modules utilized in downloading files from the web for later conversion and search operations;

5           Figure 1B is a generalized schematic drawing showing how files in various formats are converted by a utility program into Portable Document Format (PDF) files;

10           Figure 2 is a schematic flowchart showing the method in searching non-portable document format files;

            Figure 3 is a representation of a window which indicates messages to the operator for finding other matches;

15           Figure 4 is a drawing showing the basic steps involved in converting files from various different formats into PDF files and then linking them to desired portions of the original file;

20           Fig. 5 is a flow chart illustrating the conversion of a Windows Help File into Rich Text Format (RTF);

            Fig. 6 is a flow chart illustrating the conversion of HTML files to Rich Text Format (RTF);

            Fig. 7 is a flow chart showing the conversion of an HTML Help file to Rich Text Format (RTF);

25           Fig. 8 is a flow chart showing the conversion of a Rich Text Format file to Portable Document Format (PDF) files;

30           Fig. 9 is a flow chart illustrating a search which can be instituted on the PDF files after multiple file types have been converted to PDF;

            Fig. 10 is a set of selected topic files side-by-side indicating one topic file in PDF copy format and the same topic file in original copy format.



## GLOSSARY LIST

5

**10**

15

20

**Help and creates an RTF file.**

25

30



and codes used by programs and by printers or other devices. Examples of formats include RTF (Rich Text Format); DCA (Document Content Architecture); PICT, DIF (data interchange format), DXF, TIFF (tagged image file format), and EPSF (Encapsulated PostScript Format).

**FORMAT:** This involves a structure or layout of an item. Screened formats are fields on the screen; report formats are columns, headers and footers on a page. Record formats are the fields within a record. File formats are the structure of data and program files, word processing documents and graphics files (display lists and bitmaps) with all their proprietary headers and codes.

**FORMAT PROGRAM:** This is software that initializes a disk. There are two formatting levels. The low level initializes the disk surface by creating the physical tracks and storing sector identifications in them. Low level format programs lay out the sectors as required by the particular type of drive technology used (IDE, SCSI, etc.). The high-level format creates the indexes used by the operating system (Mac, DOS, etc.) to keep track of the data stored in the sectors.

**FULL-TEXT SEARCH:** Full-Text search is a mechanism for searching for text in a collection of documents using various criteria. Adobe makes this available for files released on CD-ROM and Verity for files released on Web sites. It is necessary in both these cases to create auxiliary files to support full-text search. The user may search all documents or any subset of the documents using wildcards--for example, searching for "install\*" will find all occurrences of install, installing, installation, installed, etc. The user may also use Boolean arguments--for example, searching for "installation and printers" will find all documents in which both the words "installation" and "printers" occur.

**HTM:** This is a file name extension -- for example, CONTENTS.HTM or INDEX.HTM. This extension is usually used to identify files read by an Internet browser, such as Internet Explorer or Netscape.

HTML (Hypertext Markup Language): This is a standard for defining hypertext links between documents. It is a subset of SGML (Standardized General Markup Language).

**HYPERLINK:** A hyperlink is a part of a page, whether the page is displayed from a CD-ROM or from a Web site, that the user can click with the mouse to perform some function, such as open a document, play a video, or display an external file.

awk\appl\041503L.doc

5 For example, traveling among the links to the word "iron"  
in an article might lead the user to the periodic table  
of the elements or else a map of the migration of  
metallurgy in iron age Europe. The term "hypertext" was  
coined to described documents (as presented by a  
10 computer) that expressed the non-linear structure of  
ideas as opposed to the linear format of books, films,  
and speech.

NEXT HIT OPTION: This is an option provided by a search engine to facilitate navigation from one "hit," or found item, to the next. Ordinarily, the user performs a search and the search engine presents the user with a "hit" list. This is a list of documents in which the items for which the user is searching can be found. When the user opens a document from the list, the first "hit" in the document is displayed. The user then moves to successive hits by selecting the next hit option.

awk\appl\041503L.doc

**ORIGINAL PDF:** This is a PDF file that was originally created to be delivered as a PDF file. It is usually a complete book, and it includes all graphics, special fonts, etc.

5 PDF COPY: This is a PDF file that was created from  
another type of file, such as Windows Help, HTML, or HTML  
Help. It contains only the text from the other file.

PDF FILES CREATED FROM TEXT EXTRACTED FROM OTHER FILE  
TYPES: The disclosure includes utilities that read the  
10 unformatted text from other types of files. The text is  
used to generate a PDF companion file of the original  
file that has links from each page into the corresponding  
location within the original file.

POSTSCRIPT DRIVER: This is Windows software which  
15 facilitates printing from a Windows application to a  
PostScript printer.

POSTSCRIPT FILE: This is a Windows file created by redirecting the commands generated by a PostScript driver to a file instead of to a printer. It can be copied to a PostScript printer or used by Adobe Acrobat Distiller to produce PDF files.

PREVIOUS HIT OPTION: This is an option provided by a search engine to facilitate navigation from one "hit," or found item, to the next. Ordinarily, the user performs a search and the search engine presents the user with a "hit" list. This is a list of documents in which the items for which the user is searching can be found. When the user opens a document from this list, the first "hit" in the document is displayed. The user then moves to successive hits by selecting the next hit option. Once the user has selected the next hit option, it is possible to return to the previous successive hit by selecting the previous hit option.

**RTF:** This is Rich Text Format, an adaptation of DCA (Document Content Architecture). This allows a user to transfer formatted text documents between applications, even those running on different platforms.

5    RTF FILE IN WORD:    This is the process of opening an RTF file in Word.    Word converts the RTF file into a Word document.

**RTF PAGES:** These are pages displayed in Word when it has an RTF file open. This allows the developer to see the separate pages.

SEARCH: This is the action of seeking the location of a file, or to search a file or data structure for specific data. A search is carried out by comparison or calculation to determine whether a match to some specified pattern exists or whether some other criteria have been met.

**SEARCH ALGORITHM:** This is an algorithm designed to locate a particular element, called a target in a list.

SEARCH TARGET: The search target is the text which  
20 defines what is being searched for. This could be a  
literal string of text which is to be found, such as  
"installation instructions," or a string containing  
wildcards, such as "install\*", or a string containing  
Boolean instructions, such as "installation and  
25 printers."

**SEARCH TERM:** See "Search Target."

**SENDKEYS:** This is a function supported by Visual Basic and some other programs running under Windows that permits one software application to send keystrokes to another to simulate user input.

UNFORMATTED TEXT: This term refers to text that does not contain formatting information attributes, such as font name, point size, bold, italics, underline, etc., or does not possess the structure associated with tables, columns, indented paragraphs, etc.

VERITY SEARCH ENGINE: This is a software suite developed by Verity, and used on the Unisys Support Web site, that facilitates full-text search of files on a Web site. It includes both the software that the site administrator has to execute to create files necessary to support full-text search as well as the software that the user accesses to perform the searches. Verity Inc., 894 Ross Drive, Sunnyvale, CA 94089.

WEB BROWSER: A client application that enables a user to view HTML documents on the World Wide Web, another network, or the user's computer; follow the hyperlinks among them; and transfer files. Text-based Web browsers, such as Lynx, can serve users with shell accounts but show only the text elements of an HTML document: most Web browsers, however, require a connection that can handle IP packets but will also display graphics that are in the document, play audio and video files, and execute small programs, such as Java applets or ActiveX controls, that can be embedded in HTML documents. Some Web browsers require helper applications or plug-ins to accomplish one or more of these tasks. In addition, most current Web browser permit users to send and receive e-mail and to read and respond to newsgroups.

WINDOWS: This is an operating system introduced by Microsoft Corporation in 1983. Windows is a multi-tasking graphical user interface environment that runs on MS-DOS based computers. Windows provides a standard interface based on drop-down menus, windowed regions on the screen, and a pointing device such as a mouse. The



programs used must be specially designed to take advantage of these features. A graphics-based operating system from Microsoft that provides a desktop environment similar to the Macintosh in which applications are displayed in re-sizeable moveable windows on a screen. Starting with Windows 95, the Windows system is a self-contained 32-bit operation system that requires a minimum Intel 386. In order to use all the features of Windows, applications must be written for this system.

- 10 WINDOWS HELP: Windows-based help systems are automated Windows utilities that provide procedural and system information to software users in lieu of paper-based documentation. Windows-based help supports context-sensitive help, which lets the user access topics in a help file that are relevant to the user's location in the application.

DESCRIPTION OF PREFERRED EMBODIMENT:

Fig. 1A is a generalized drawing which illustrates the environmental modules which constitute the operating modules which permit the conversion of downloaded multiple-type files from the Web into Portable Document Format (PDF) files for observation on a observable window by the operator.

Now referring to Fig. 1A, a personal computer 10 is seen having a memory 12 and operating system 14 and is also connected to a disk storage unit 16.

The personal computer 10 (user workstation) is provided with an Adobe Acrobat program 22.

The World Wide Web 5 is seen connected to the personal computer 10 and may download digital data in various different formats.

A Verity Search Engine 9 connected to the terminal server 8 can initiate a search on the Web 5 and bring about a download of multiple files to the user workstation 10. However, some of these files may be in one particular format, while others may be in different formats, thus instigating a problem when a browser or search engine is used in order to find a particular subject matter or topic on any one of the particular files.

Fig. 1B is an overall generalized drawing showing the basic steps in the creation of text copies from various types of downloaded files for conversion into Portable Document Format, or PDF files. For example, as seen in Fig. 1A, the Windows Help file (W1) is converted by a utility program (U2) into a Portable Document Format copy designated (WC).

Again, in Fig. 1A, a hypertext mark-up language file (HTML) designated as (M1) is passed through a utility program (U2M) after which there is provided at

Further, in Fig. 1A, there is seen an HTML Help file (HH1) which is passed through a utility program (U2HH) in order to provide a Portable Document Format copy designated (HHC).

Now referring to Fig. 2, there is seen a generalized view for the searching of non-Portable Document Format files. Here, it is desired that a search be made on a particular topic or target such as "I/O" for example, in order to finally provide and display the data of the original file on that particular topic. Thus, as seen in Fig. 2, at step (NP1), there is instituted a search of all of the Portable Document Format (PDF) files.

At step (NP3), the operator can click a button which appears on that particular page that is displayed, and then at step (NP4), the operator can open the original file to the selected topic, for example, such that the original target topic, such as "I/O" will now be displayed and seen in its original file form.

Seen on this window is a set of icons, one of which can be pressed for "search" and another icon which  
35 can be pressed for search results. Then, there is another icon which shows a way to find the previous match

The search results icon will provide a display of a list of documents that contain matches, while the search icon is used to change the search topics.

A sequence of original files are shown which are to be the object of a search. The Windows Help files are designated W1 and the HTML files are designated M1, while the HTML Help files are designated HH1, and the Help file is designated H1.

The next level of steps shown respectively, as W3, M3, HH3, and H3, all involve the step of conversion with use of the Adobe Acrobat software converter.

Then in Fig. 4, there is seen step W5 which involves two separate functions, one of which is the set of buffers to hold the PDF files, together with an explanation message regarding the files in the buffer. An example of an explanation message and a link created by this program are shown in the left panel of Figure 10.

awk\appl\041503L.doc

As will be seen in the next succeeding set of drawings, it should be understood that there are certain intermediate steps involved, whereby the original files are first converted to Rich Text Format (RTF), after which the subsequent RTF files can then later be converted to Portable Document Format (PDF).

15           At step W2, the program will open the Windows  
Help file.

At step W4, the program will then go to the list to read the number of the next topic that has a Topic ID. For example, this next topic might be the subject of "Channel Adapters".

awk\appl\041503L.doc

Then at step W7, the program will copy the text from the Clipboard and format the Rich Text Format pages, after which there is a return to step W4 in order to get the text from the next topic.

5           Fig. 6 is a flow chart illustrating the steps involved for converting the HTML files to Rich Text Format (RTF). At step 1, the program will acquire the name of the directory containing the HTML files and also the name of the Output Rich Text Format (RTF) file. Note  
10 that an HTML "document" can consist of a number of files with the HTM extension.

          Then at step M2, the program will get the next file in the directory with the HTM extension. This is a Windows/DOS file name extension, which is equivalent to  
15 HTM, as for example, CONTENTS.HTM or INDEX.HTM. This extension is usually used to identify files read by an Internet browser, such as Internet Explorer or by Netscape.

          At step M3, a decision block is presented which  
20 presents the query as to whether or not another file with the HTM extension is present. If the answer is (NO), then the program will end at step M3E. If the answer is (YES) at step M3, then step M4 occurs to open the particular file with the ActiveX control which will use  
25 the InnerText method to read the text. InnerText is a software mechanism within the Microsoft ActiveX control that supports Internet Explorer and will extract unformatted text from within the body of a HTML file.

          Then, at step M5, the program will format the  
30 Text into Rich Text Format pages (RTF).

          After step M5, the program loops back to step M2 to get the next file in the directory with the HTM extension.

          Fig. 7 is a flow chart illustrating the  
35 conversion of an HTML Help file into a Rich Text Format (RTF) file. An HTML Help file is also called a CHM file

Here at step HH1, the program will acquire names of the CHM file directory, which contains the HTML files from which the CHM file is constructed and the Output RTF file to be created by the program.

At step HH3, a query block is presented to query whether an additional file with an HTM extension is present. If the answer is (NO), then the program ends here at step HHE. If the answer is (YES), that is to say, a file is present, then at step HH4, the program will open the file with the ActiveX control and use the InnerText method to read the text. This copies unformatted text from within the body of a HTML file. Graphics, font information, such as point size, bold, italic, etc., and structure, such as tables, columns, etc., are not copied.

25           After this, the program loops from HH5 back to  
HH2 in order to operate on the next file in the  
directory.

At step CRP1, the program will open the Rich  
35 Text Format file in Word so that the Word program of

At step CRP2, the program will use the Word program to print to file, using a PostScript driver. The PostScript driver is a portion of Windows software which facilitates printing from a Windows application to a PostScript printer.

15           At step CRP4, the program will open the  
PostScript file in the Adobe Acrobat Distiller.

20 With the development of the PDF file as shown  
in Fig. 8, the Portable Document File can now relate to  
Fig. 4 which shows the level of Portable Document Format  
files seen at steps W4, M4, HH4, and H4.

This can further be expounded by the flow chart seen in Fig. 9, where now that the Portable Document Format (PDF) copies have now been isolated, then a search can be initiated using the Adobe Acrobat programs.



Now referring to Fig. 9 at step S1, the program will initiate a search of a particular topic through the Adobe Acrobat program.

Then at step S2, there is presented a list of the Portable Document Format (PDF) documents, showing the list of hits to the user.

At step S3, the user selects a Portable Document Format document and opens it to the first hit.

At step S4, a decision box is initiated to query of whether the file is originally a Portable Document File. If the answer is (YES), then the program sequence is to step S7 to query whether the search should end.

At step S4, if the answer is (NO), that is to say, the file is not originally a Portable Document Format file, then at step S5 the user will click the "Open Document" button on the top of the display page.

At step S6, the original document is now opened to the particular topic containing the text in the Portable Document Format file.

At step S7, a decision box presents the question of whether this is the end of the search. If the answer is (YES), the search ends at step S7E. If it is not the end of the search (NO), then step 8 occurs where the user clicks the "next hit" button on the tool bar of the Portable Document Format file.

Then, step S8 loops back to step S4 in order to continue through S5, S6 and S7 until the search has ended at S7E.

Now referring to Fig. 10, there is illustrated a page of unformatted text which is shown on the left side of the page, and its corresponding original file which is indicated on the right-hand side of the page.

As an example, the subject matter was that of "Establishing a named pipe to a COMs Application". Here, it will be noticed that the unformatted text does not

contain all the information, such as graphics, etc., but that the original file shown on the right-hand side shows the original text together with the graphics and detailed material which may not appear in the unformatted text.

5           Thus, it can now be understood that a series of document information such as articles, books or manuals can be downloaded from the Web and exist in different types of formats. This normally would make it unwieldy or impossible to search through the entire list of  
10 downloaded documents in order to get information on a particular topic that was desired since any one particular search browser is specific to the handling of any one particular format, but not available or useful in handling the many different format types involved, or  
15 multiple types of formats.

          Thus, the present system, by using the intermediate step of providing the Rich Text Format which can then be converted to the Portable Document Format, and then the Portable Document Format is utilized as  
20 being compatible with and accessible to search purposes by use of the Adobe Acrobat program, the multiple numbers of different files, documents, articles or pages downloaded from the Web via the Verity Search Engine can now be searched for a given topic and then displayed in  
25 Portable Document Format (PDF).

          Then subsequently, the Portable Document Format (PDF) can then be linked back to the original text of the original pages holding the desired topic information desired by the user and these can be displayed in their  
30 original format with full graphics, colors, lists, tables and any other types of display which would not be available in the PDF format.

          While a particular implementation of the above-described invention has been shown in a particular  
35 effective implementation, there may be other implementations of the invention which are derivable from

[illegible]